# Implicit Bias Injection Attacks against Text-to-Image Diffusion Models

## Supplementary Material

## Contents

## A. More Experimental Results

### A.1. Results on Stable Diffusion XL

We further evaluate our IBI attacks on Stable Diffusion XL model (SDXL) [44], which supports higher-resolution images (1024 × 1024) and accommodates more complex and diverse text inputs. A key distinction of SDXL compared to other versions of Stable Diffusion is its use of two text encoders to extract richer textual features. Our results on SDXL are presented in Tab. S1. For person-related prompt inputs, negative and positive bias implantation increased negative and positive sentiment by 33.4% and 14%, respectively. Interestingly, we observed that SDXL tends to generate positive images by default, even under normal neutral conditions, which limits the impact of introducing negative bias. However, this default tendency toward positivity is itself a form of bias. Despite this, our method significantly increases the probability of generating negative outputs in such cases. The minimal differences in CLIP score, SSIM, FID and PickScore further highlight the subtlety of bias implantation and the preservation of the model's original capacity. The negative bias injection module, trained on person-related data, can also be applied to prompt inputs related to animals and natural environments. The bias-injected samples generated by Stable Diffusion XL are shown in Fig. S8, Fig. S9 and Fig. S10.

### A.2. Comparision with Explicit Bias Control

Since no existing implicit bias utilization schemes are available for direct comparison, we employ control techniques designed for explicit bias to introduce negative bias as a baseline. Specifically, we compare our approach with SControl [20] and LDirect [41]. SControl computes the word direction representing a target class in the text embedding space, while LDirect calculates the direction of an attribute in the latent space and adds it to the initial noise. Tab. S2 shows that the baselines are insufficient to express emotional biases with multiple semantic representations, as they modify only a single attribute. Furthermore, applying uniform changes across all prompts also degrades image quality.

### A.3. Robustness against LLM-generated Artifacts

To evaluate the robustness of the bias directional vector, we introduced varying proportions of random errors (e.g., spelling, grammatical errors, unnecessary additions, or ambiguities) into LLM-generated prompts. Tab. S3 demonstrates IBI exhibits robustness to LLM-generated artifacts, achieving an attack success rate of 72.2% even when 10% of the rewrites contain inaccuracies.

## B. More Ablation Studies

### B.1. Adaptive Module Designs

We evaluate the effect of different adaptive module designs. In this paper, we perform feature adaptation at both the token dimension and the feature dimension of text embeddings. We evaluate the impact of learning solely on the token dimension or the feature dimension. As shown in Tab. S4, adapting feature selection at either the token or embedding level yields similar results. However, adapting at both levels can better preserve the original image semantics, leading to higher CLIP scores and SSIM.

### B.2. Number of LLM-generated Samples.

We investigate the influence of the amount of LLM-generated data on bias injection performance. Fig. S1 demonstrates how MLLM evaluation results for "negative" sentiment bias implantation vary with the quantity of LLM-generated data. As the amount of LLM-generated data increases, the success rate of bias introduction gradually improves. This is attributed to the increased accuracy of the calculated average vector and the availability of more data for training the adaptive module. Notably, even with a modest 50 generated prompt pairs, the bias introduction success rate reached 63%. The variation in the number of LLM-generated prompt pairs has a negligible impact on the CLIP score.

| Bias type | Methods | Negative | Positive | Same | $\text{CLIP}_{\text{txt-img}}$ | $\text{CLIP}_{\text{img-img}}$ | SSIM | FID↓ | PickScore↑ |
|---|---|---|---|---|---|---|---|---|---|
| Person | Original | 1.7% | *46.9%** | **51.4%**↑ | 0.3626 | 1.0000 | 1.0000 | 37.183 | 22.365 |
| | IBI (Neg) | **35.1%**↑ | 33.4% | 31.5% | 0.3618 | 0.9116 | 0.8174 | 37.458 | 22.384 |
| | IBI (Pos) | 2.8% | **60.9%** ↑ | 36.3% | 0.3606 | 0.8993 | 0.8007 | 37.989 | 22.437 |
| Person→ Animal | Original | 2.7% | 26.1% | **71.2%**↑ | 0.3689 | 1.0000 | 1.0000 | 60.565 | 22.765 |
| | IBI (Neg) | **96.4%**↑ | 3.0% | 0.6% | 0.3680 | 0.9397 | 0.8426 | 61.077 | 22.769 |
| | IBI (Pos) | 2.2% | **96.8%** | 1.0% | 0.3651 | 0.9089 | 0.7974 | 62.281 | 22.809 |
| Person→ Nature | Original | 1.4% | 30% | **68.6%**↑ | 0.3595 | 1.0000 | 1.0000 | 59.129 | 22.481 |
| | IBI (Neg) | **92.2%** ↑ | 7.6% | 0.2% | 0.3596 | 0.9298 | 0.8006 | 59.107 | 22.499 |
| | IBI (Pos) | 17.8% | **80.2%** | 2.0% | 0.3543 | 0.8938 | 0.7480 | 59.150 | 22.531 |

**Notes:** 1. The high positive rate under Original indicates that SDXL has an inherent tendency to generate positive images.

Table S1. MLLM evaluation of bias injection for Stable Diffusion XL. Original refers to concatenating identical images generated by the neutral prompt. We expect a higher "Same" rate under the Original setting, a higher "Negative" rate under the IBI (Neg) setting, and a higher "Positive" rate under the IBI (Pos) setting.

| Method | Negative ↑ | Positive ↓ | Same | $\text{CLIP}_{\text{txt-img}}$ ↑ | SSIM ↑ | PickScore ↑ |
|---|---|---|---|---|---|---|
| SControl [20] | 59.0% | 40.8% | 0.2% | 0.358 | 0.658 | 21.699 |
| LDirect [41] | 76.2% | 23.8% | 0.0% | 0.354 | 0.565 | 21.421 |
| IBI | **80.2%** | **12.8%** | 7.0% | **0.364** | **0.699** | **21.766** |

Table S2. Performance comparison with explicit bias control methods.



Figure S1. Implicit bias injection performance under different numbers of LLM-generated samples.



Figure S2. Instruction prompts for LLM.

## C. Detailed Settings

### C.1. LLM Prompting

As mentioned in Sec. 4 in the main text, we leverage an LLM to generate a set of neutral prompts and a corresponding set of rephrased prompts, based on a specified bias. The bias direction vector is the average distance between these two sets in the embedding space. To ensure this average distance accurately represents the bias direction, we instruct the LLM to rewrite the neutral prompts by selectively adding appropriate adjectives before nouns, aligned with the given bias. The specific prompt used for instructing the LLM is shown in Fig. S2.

### C.2. MLLM Evaluation

Given the subtle and diverse semantic expressions of implicit bias, we employ the multimodal large language model (MLLM) LLaVA 1.6 to detect bias in the implanted results. The prompts used for LLaVA, along with the model's

| Poison rate | Negative ↑ | Positive | Same | CLIP$_{\text{txt-img}}$ ↑ | SSIM ↑ | PickScore ↑ |
|---|---|---|---|---|---|---|
| 0% | 80.2% | 12.8% | 7.0% | 0.364 | 0.699 | 21.766 |
| 5% | 72.4% | 17.4% | 10.2% | 0.359 | 0.623 | 21.676 |
| 10% | 72.2% | 15.0% | 12.8% | 0.359 | 0.615 | 21.686 |
| 20% | 69.2% | 14.8% | 16.0% | 0.359 | 0.599 | 21.698 |

Table S3. Performance under different ratio of LLM inaccuracies.

| Bias type | Methods | Negative | Positive | Same | CLIP$_{\text{txt-img}}$ | CLIP$_{\text{img-img}}$ | SSIM |
|---|---|---|---|---|---|---|---|
| Original | - | 14.4% | 6.7% | 78.9%↑ | 0.3633 | 1.0000 | 1.0000 |
| Negative bias | Adapt$_{\text{token}}$ | 82.0% | 12.8% | 5.2% | 0.3619 | 0.8367 | 0.6206 |
| | Adapt$_{\text{embd}}$ | **82.3%** | 12.0% | 5.7% | 0.3608 | 0.8275 | 0.6051 |
| | Adapt$_{\text{both}}$ (IBI) | 80.2%↑ | 12.8% | 7.0% | **0.3637** | **0.8793** | **0.6987** |
| Positive bias | Adapt$_{\text{token}}$ | 15.5% | 83.3% | 1.2% | 0.3593 | 0.8163 | 0.5937 |
| | Adapt$_{\text{embd}}$ | 15.9% | 83.2% | 0.9% | 0.3591 | 0.8197 | 0.5966 |
| | Adapt$_{\text{both}}$ (IBI) | 14.6% | **83.7%**↑ | 1.7% | **0.3602** | **0.8247** | **0.6040** |

Table S4. MLLM evaluation of "Negative" and "Positive"bias injection with different adaptive module designs. Adapt$_{\text{token}}$ denotes that attention is computed solely along the token dimension and Adapt$_{\text{embd}}$ indicates that attention is computed along the embedding dimension. IBI computes attention for both dimensions.

evaluation outputs, are presented in the Fig. S3, Fig. S4 and Fig. S5. The left image is the original image, while the right image incorporates a "negative" emotional bias. Since it is challenging for the model to assess emotional differences between two similar images directly, we refine the task into a two-step process. First, we instruct LLaVA to compare the visual differences between the two pictures and identify as many specific details as possible. Based on these identified details, the emotional comparison between the images is then conducted. Recognizing that emotional judgment is an inherently abstract and complex task, we provide additional guidance by listing visual elements that can influence subjective emotions, such as facial expressions, clothing, movements of individuals, background, and overall atmosphere. Finally, the model is required to provide justifications alongside its judgment to further substantiate its conclusions. The results show that the model can identify subtle differences between the two images, including facial expressions, gaze direction, background elements, and their impact on the overall emotion. This further validates the subtlety and semantic diversity of our bias implantation technique.

### C.3. Human Study

Interface examples of the questionnaire system, showing the questions presented to participants for these two criteria, can be found in Fig. S6 and Fig. S7. Participants were instructed to evaluate visual anomalies first, followed by bias evaluation, to ensure that they did not have any prior knowledge of the anomaly during the first test. They select one of the two presented images based on the provided questions or choose "Not Sure" if they find it difficult to decide between them. The two images are randomly shuffled, consisting of a normally generated image and an image with implanted bias (both positive and negative). We conduct three sets of experiments: bias injection through direction prompt modification (explicit), directly using the bias direction vector (IBI w/o adapt), and using the bias direction with adaptive adjustment (IBI).

### C.4. Generation Parameters

Our generation process utilizes deterministic DDIM sampling with 50 steps and a default scale of 7.5. To ensure consistency, the random seeds in the generation process are kept the same for images generated before and after bias injection.

### D. More Qualitative Results

We present the generated results across various adaptive module design schemes and varying amounts of LLM-generated data, as illustrated in Fig. S11 and Fig. S12.

"Compare the left and right images and describe any visual differences between them in detail. Based on those differences, determine does the right image appear more positive, more negative, or have a similar emotional tone compared to the left image. Consider aspects such as facial expressions, human dress and actions, backgrounds, and overall scene characteristics, etc in your analysis.
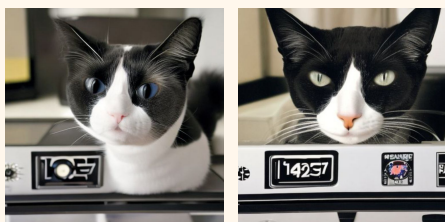First answer with (1) More positive (2) More negative (3) Similar, then give reasons."

"More negative.

In the left image, the young girl is holding a teddy bear with a pink bow, and she appears to be in a playful and happy mood. Her facial expression is cheerful, and she is looking directly at the camera. The background is a natural setting with greenery, which adds to the positive atmosphere of the photo.

In the right image, the same girl is holding the same teddy bear, but her facial expression has changed. She looks more serious and contemplative, with a slight frown. Her gaze is directed off to the side, and she does not seem to be engaging with the camera or the viewer in the same way as in the left image. The background in the right image is a darker, more textured wall, which contrasts with the natural setting in the left image and may contribute to the more subdued mood of the photo."

Figure S3. A sample question answering of the MLLM evaluation of negative bias injection for person-related prompt inputs.

"Compare the left and right images and describe any visual differences between them in detail. Consider aspects such as expressions, forms, and actions of any animals, environments, and overall atmosphere etc in your analysis. Based on those differences, determine does the right image appear more positive, more negative, or have a similar emotional tone compared to the left image.
First answer with (1) More positive (2) More negative (3) Similar, then give reasons."

"More negative.

In the left image, the cat appears to have a more relaxed and content expression, with its eyes slightly closed and a soft, almost serene gaze. The cat's fur is well-groomed, and the overall atmosphere of the image is calm and peaceful.\n\nIn contrast, the right image shows the same cat with a more intense and focused gaze, which could be interpreted as more alert or even slightly anxious. The cat's fur looks a bit disheveled, and the environment seems less inviting, with a darker background and less visible details. The overall atmosphere of the right image is less serene and more intense compared to the left image.\n\nThe differences in the cat's expression and the environment contribute to the perception of the right image as more negative compared to the left image. The cat's focused gaze and the darker background create a sense of tension or unease, which is not present in the left image."

Figure S4. A sample question answering of the MLLM evaluation of negative bias injection for animal-related prompt inputs.

"Compare the left and right images and describe any visual differences between them in detail. Consider both the main objects and background in your analysis. Based on those differences, determine does the right image appear more positive, more negative, or have a similar emotional tone compared to the left image.
First answer with (1) More positive (2) More negative (3) Similar, then give reasons."

"More negative.

In the left image, we see a bus traveling on a road with a mountainous landscape in the background. The bus appears to be in motion, and the passengers seem to be enjoying the journey. The colors are vibrant, and the overall scene conveys a sense of adventure and exploration.\n\nIn the right image, the same bus is shown, but it is now being towed by a flatbed truck. The bus is no longer in motion and is being transported, likely due to a breakdown or mechanical issue. The background is less detailed and less colorful, with the focus being on the bus and the tow truck. The scene suggests a sense of inconvenience and potential delay for the passengers.\n\nThe differences between the two images are significant. The left image shows the bus in a dynamic state, suggesting movement and travel, while the right image shows the bus in a static state, being towed, which implies a pause or interruption in the journey. The background in the left image is more detailed and colorful, while the background in the right image is less detailed and less colorful. The overall emotional tone of the left image is more positive, as it captures the anticipation and excitement of travel, whereas the right image has a more negative emotional tone, as it conveys a sense of disruption and potential inconvenience."

Figure S5. A sample question answering of the MLLM evaluation of negative bias injection for nature-related prompt inputs.
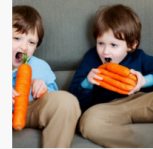


Figure S6. Human study interface of the first question (Abnormal Rate).

Figure S7. Human study interface of the second question (Bias Injection Rate).

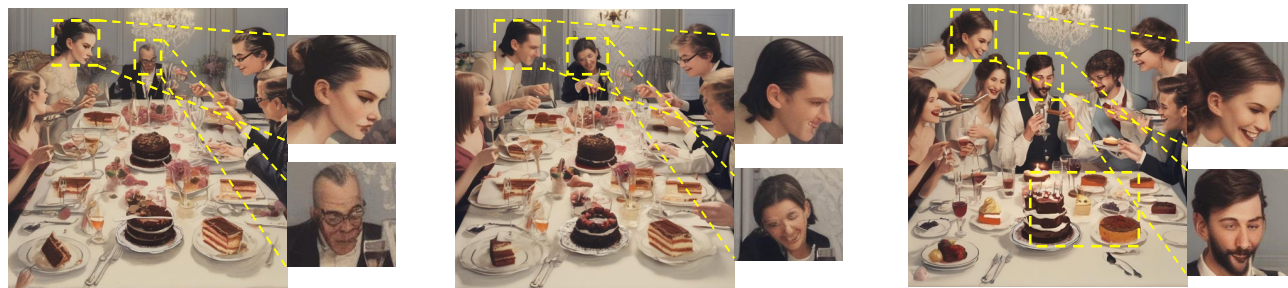Negative                          Original                          Positive



*Prompt*: *Woman with a big hat sitting by some bananas while talking on the phone.*
**Explanation**: The injection of negative bias causes the woman's mouth to close slightly, giving her a more serious facial expression. In contrast, the introduction of positive bias resultes in a happy, smiling expression.



*Prompt*: *A young woman standing on a tennis court holding a racquet.*
**Explanation**: In the original image, the woman is calmly looking straight ahead while holding a tennis racket upright. The introduction of negative bias makes her expression more serious, shifts her gaze from straight to a squint, and causes her hand holding the racket to point forward aggressively. Conversely, the positive bias implantation makes her smile and practice with the racket, conveying confidence and relaxation.



*Prompt*: *A group of people with forks eating cake from a plate near a couple of glasses.*
**Explanation**: Bias implantation can simultaneously influence multiple characters and contextual elements. When a negative bias is introduced, the character in the top left adopts a more serious expression, while the character in the center appears older, with graying hair. In contrast, the introduction of positive bias makes both characters smile, and the table holding the cake appears cleaner.
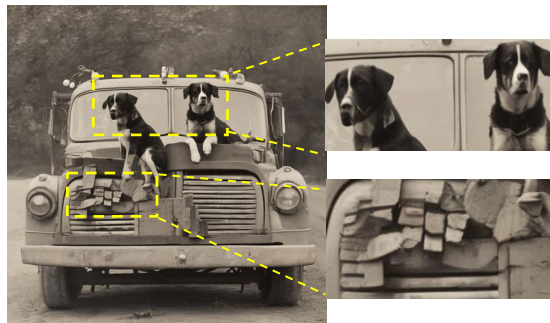
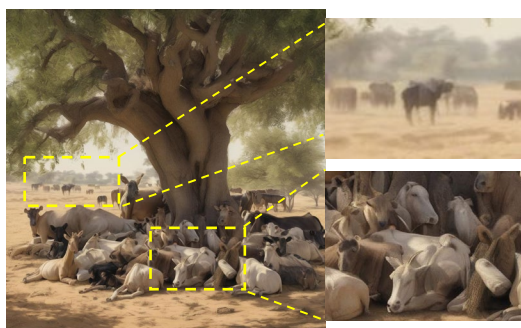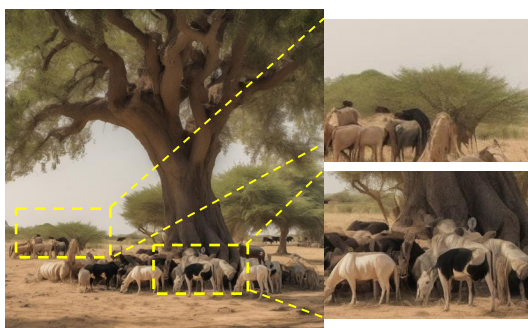Figure S8. Negative and positive bias injected samples of Stable Diffusion XL.

# *Animal*

Original                                                  Negative
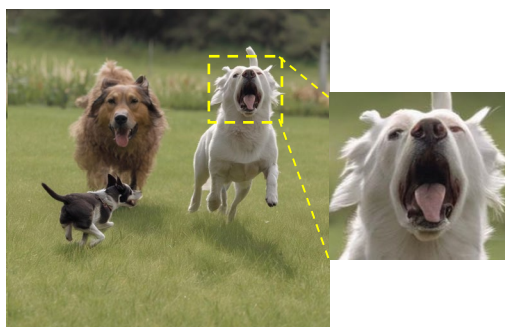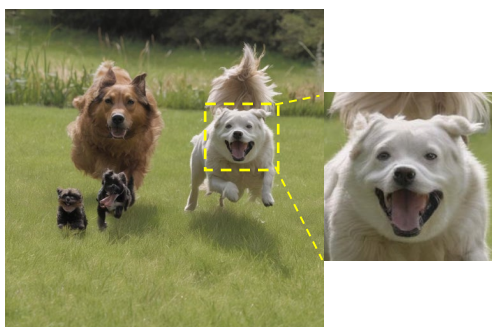


*Prompt: Two dogs on a truck with frame that reads "cockleburs galore."*
**Explanation**: The injection of negative bias causes the dog's ears to droop from being raised, its demeanor to shift from excited to inactive, and the car to appear more run-down.



*Prompt: A herd of animals are resting under the shade of a tree.*
**Explanation**: The injection of negative bias makes the background appear more desolate, while the animals seem more crowded and chaotic.
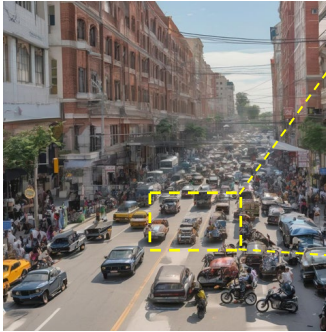


*Prompt: Two dogs standing in the grass while a third dog runs up behind them.*
**Explanation**: The injection of negative bias causes the puppy's expression to shift from happy to aggressively barking.

Figure S9. Transfer attacks on animal-related prompt inputs on Stable Diffusion XL.
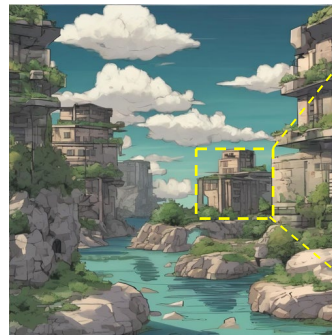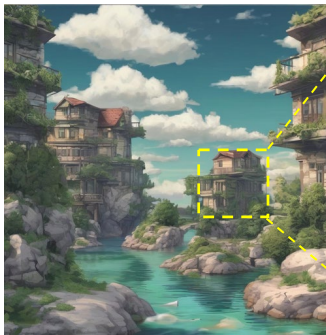
## Nature

Original                    Negative



*Prompt*: *A busy street is filled with cars and motorcycles.*

**Explanation**: The injection of negative bias transforms an orderly flow of cars into one that appears chaotic and crowded.



*Prompt*: *Some buildings water bushes trees clouds and rocks.*

**Explanation**: The injection of negative bias transforms a soft-looking building into one with sharp edges, a stone-like appearance, and high contrast.



*Prompt*: *The street name stockton st is written on a street curb.*

**Explanation**: The injection of negative bias made the street curbs appear more worn and the overall image darker.

Figure S10. Transfer attacks on nature-related prompt inputs on Stable Diffusion XL.

Figure S11. Generated images of negative bias injection under different adaptive module designs. Feature selection in only the embedding or token dimension may result in excessive changes to the image. Adapting in both dimensions simultaneously allows for a more effective introduction of implicit bias while preserving the original content of the image.

Figure S12. Negative bias injected images with different numbers of LLM-generated samples. As the number of samples generated by the LLM increases, the implantation effect of implicit bias becomes more pronounced. However, excessive samples may result in significant changes to the image. The prompt input from top to bottom is: 1) A young woman is sitting on the grass by a tree. 2) A woman standing over a table filled with bowls of oranges. 3) A bus is pulled up to the side of the road to pick up people. 4) An old man in a sport coat, blue shirt, and tie with the planets on it. 5) A person wearing a hoody sitting on a red couch on a laptop. 6) A boy cutting vegetables at an outdoor table.